

PANDAS FOR BEGINNERS

AIOTLAB, September 2023



Table of Contents

- 1 Getting Started
- 2 Introduction to pandas Data Structures
- 3 Essential Functionality
- 4 Summarizing and Computing Descriptive Statistics



Table of Contents

1 Getting Started

2 Introduction to pandas Data Structures

3 Essential Functionality

4 Summarizing and Computing Descriptive Statistics



Getting Started

- **pandas** contains **data structures** and **data manipulation tools** designed to make *data cleaning* and *analysis* fast and convenient in Python
- **pandas** is often used in tandem with:
 - Numerical computing tools (**NumPy** and **SciPy**)
 - Analytical libraries like (**statsmodels** and **scikit-learn**)
 - Data visualization libraries (**matplotlib**)
- **pandas** adopts significant parts of NumPy's idiomatic style of array-based computing, especially array-based functions and a preference for data processing *without* **for** loops
- While **pandas** adopts many coding idioms from NumPy, the biggest difference is that pandas is designed for working with **tabular** or **heterogeneous data**
 - **NumPy**, by contrast, is *best* suited for working with **homogeneously** typed **numerical** array data



Getting Started

- The following *import conventions* for **NumPy** and **pandas** are commonly used:

```
1 import numpy as np  
2  
3 import pandas as pd
```

- It is easier to import Series and DataFrame into the local namespace since they are so frequently used

```
1 from pandas import Series, DataFrame
```



Table of Contents

1 Getting Started

2 Introduction to pandas Data Structures

3 Essential Functionality

4 Summarizing and Computing Descriptive Statistics



Series

- A **Series** is a **one-dimensional** array-like object containing a sequence of *values* (of similar types to NumPy types) of the **same type** and an associated array of *data labels*, called its **index**

```
1  obj = pd.Series([4, 7, -5, 3])
2
3  obj
4
5  Output:
6  1    7
7  2   -5
8  3    3
9  dtype: int64
```



Series

- The string representation of a Series displayed interactively shows the index on the left and the values on the right
 - Since we did not specify an index for the data, a default one consisting of the integers **0** through **N - 1** (where **N** is the length of the data) is created
 - Can get the array representation and index object of the Series via its **array** and **index** attributes, respectively
 - The result of the **.array** attribute is a **PandasArray** which usually wraps a NumPy array but can also contain special extension array types

```
1 obj.array
2 Output:
3 <NumpyExtensionArray>
4 [4, 7, -5, 3]
5 Length: 4, dtype: int64
6
7 obj.index
8 Output: RangeIndex(start=0, stop=4, step=1)
```



Series

- Often, we'll want to create a Series with an **index** identifying each data point with a **label**

```
1 obj2 = pd.Series([4, 7, -5, 3], index=["d", "b", "a", "c"])
2
3 obj2
4 Output:
5 d    4
6 b    7
7 a   -5
8 c    3
9 dtype: int64
10
11 obj2.index
12 Output: Index(['d', 'b', 'a', 'c'], dtype='object')
```



Series

- Compared with NumPy arrays, we can use labels in the index when selecting *single* values or a *set* of values

```
1 obj2["a"]
2 Output: -5
3
4 obj2["d"] = 6
5
6 obj2[["c", "a", "d"]]
7 Output:
8 c    3
9 a   -5
10 d    6
11 dtype: int64
```

- Here `["c", "a", "d"]` is interpreted as a *list of indices*, even though it contains strings instead of integers



Series

- Using NumPy functions or NumPy-like operations, such as filtering with a Boolean array, scalar multiplication, or applying math functions, will preserve the index-value link

```
1 obj2[obj2 > 0]
2 Output:
3 d      6
4 b      7
5 c      3
6 dtype: int64
7
8 obj2 * 2
9 Output:
10 d     12
11 b     14
12 a    -10
13 c      6
14 dtype: int64
```

```
1 np.exp(obj2)
2 Output:
3 d      403.428793
4 b     1096.633158
5 a      0.006738
6 c     20.085537
7 dtype: float64
```



Series

- Can think of Series as a *fixed-length, ordered dictionary*, as it is a mapping of index values to data values => Can be used in many contexts where we might use a dictionary

```
1 "b" in obj2
2 Output: True
3
4 "e" in obj2
5 Output: False
```



Series

- Should we have data contained in a Python dictionary, we can create a Series from it by passing the dictionary
- A Series can be converted back to a dictionary with its **to_dict** method

```
1 sdata = {"Ohio": 35000, "Texas": 71000, "Oregon": 16000, "Utah": 5000}
2
3 obj3 = pd.Series(sdata)
4
5 obj3
6 Output:
7 Ohio      35000
8 Texas     71000
9 Oregon    16000
10 Utah      5000
11 dtype: int64
12
13 obj3.to_dict()
14 Output: {'Ohio': 35000, 'Texas': 71000, 'Oregon': 16000, 'Utah': 5000}
```

Series

- When only passing a dictionary, the index in the resulting Series will respect the order of the keys according to the dictionary's **keys** method, which depends on the key insertion order
 - Can override this by passing an **index** with the dictionary keys in the order we want them to appear in the resulting Series



Series

- Here, three values found in **sdata** were placed in the appropriate locations
 - But since no value for "**California**" was found, it appears as **NaN** (Not a Number), which is considered in pandas to mark missing or *NA* values
 - Since "**Utah**" was not included in **states**, it is excluded from the resulting object

```
1 sdata = {"Ohio": 35000, "Texas": 71000, "Oregon": 16000, "Utah": 5000}
2 states = ["California", "Ohio", "Oregon", "Texas"]
3
4 obj4 = pd.Series(sdata, index=states)
5
6 obj4
7 Output:
8 California      NaN
9 Ohio          35000.0
10 Oregon        16000.0
11 Texas         71000.0
12 dtype: float64
```

Series

- The **isna** and **notna** functions in pandas should be used to detect missing data
 - Series also has these as **instance methods**

```
1 pd.isna(obj4)
2 Output:
3 California      True
4 Ohio            False
5 Oregon          False
6 Texas           False
7 dtype: bool
8
9 pd.notna(obj4)
10 Output:
11 California     False
12 Ohio            True
13 Oregon          True
14 Texas           True
15 dtype: bool
```

```
1 obj4.isna()
2 Output:
3 California      True
4 Ohio            False
5 Oregon          False
6 Texas           False
7 dtype: bool
8
9 obj4.notna()
10 Output:
11 California     False
12 Ohio            True
13 Oregon          True
14 Texas           True
15 dtype: bool
```



Series

- A useful Series feature for many applications is that it **automatically aligns** by **index label** in arithmetic operations
 - This feature is similar to a *join operation* in databases

```
1 obj3
2 Output:
3 Ohio      35000
4 Texas     71000
5 Oregon    16000
6 Utah      5000
7 dtype: int64
8
9 obj4
10 Output:
11 California      NaN
12 Ohio            35000.0
13 Oregon          16000.0
14 Texas           71000.0
15 dtype: float64
```

```
1 obj3 + obj4
2 Output:
3 California      NaN
4 Ohio            70000.0
5 Oregon          32000.0
6 Texas           142000.0
7 Utah            NaN
8 dtype: float64
```



Series

- Both the Series **object** itself and its **index** have a **name** attribute
 - It integrates with other areas of pandas functionality

```
1 obj4.name = "population"
2
3 obj4.index.name = "state"
4
5 obj4
6 Output:
7 state
8 California      NaN
9 Ohio            35000.0
10 Oregon          16000.0
11 Texas           71000.0
12 Name: population, dtype: float64
```



Series

- A Series's index can be altered in place by assignment

```
1  obj
2  Output:
3  0    4
4  1    7
5  2   -5
6  3    3
7  dtype: int64
8
9  obj.index = ["Bob", "Steve", "Jeff", "Ryan"]
10
11 obj
12 Output:
13 Bob      4
14 Steve    7
15 Jeff     -5
16 Ryan     3
17 dtype: int64
```



DataFrame

- A **DataFrame** represents a **rectangular table** of data and contains an **ordered, named** collection of columns, each of which can be a **different value type** (numeric, string, Boolean, etc.)
 - The DataFrame has both a row and column index
 - Can be thought of as a dictionary of Series all sharing the same index
- **Note:** While DataFrame is physically two-dimensional, we can use it to represent *higher dimensional* data in a tabular format using:
 - Hierarchical indexing,
 - and an ingredient in some of the more advanced data-handling features in pandas



DataFrame

- There are many ways to **construct** a DataFrame, though one of the **most common** is from a **dictionary** of equal-length lists or NumPy arrays
 - The resulting DataFrame will have its index assigned automatically, as with Series
 - The columns are placed according to the order of the keys in **data** (depends on their insertion order in the dictionary)

```
1 data = {"state": ["Ohio", "Ohio", "Ohio", "Nevada", "Nevada", "Nevada"],  
2      "year": [2000, 2001, 2002, 2001, 2002, 2003],  
3      "pop": [1.5, 1.7, 3.6, 2.4, 2.9, 3.2]}  
4 frame = pd.DataFrame(data)  
5  
6 frame  
7 Output:  
8      state  year  pop  
9  0      Ohio  2000  1.5  
10 1      Ohio  2001  1.7  
11 2      Ohio  2002  3.6  
12 3    Nevada  2001  2.4  
13 4    Nevada  2002  2.9  
14 5    Nevada  2003  3.2
```

DataFrame

- For large DataFrames:
 - The **head** method selects only the **first 5 rows**
 - The **tail** method selects only the **last 5 rows**

```
1 frame.head()  
2 Output:  
3      state  year  pop  
4  0    Ohio  2000  1.5  
5  1    Ohio  2001  1.7  
6  2    Ohio  2002  3.6  
7  3  Nevada  2001  2.4  
8  4  Nevada  2002  2.9
```

```
1 frame.tail()  
2 Output:  
3      state  year  pop  
4  1    Ohio  2001  1.7  
5  2    Ohio  2002  3.6  
6  3  Nevada  2001  2.4  
7  4  Nevada  2002  2.9  
8  5  Nevada  2003  3.2
```



DataFrame

- If we specify a sequence of columns, the DataFrame's columns will be arranged in that order

```
1 pd.DataFrame(data, columns=["year", "state", "pop"])
2 Output:
3      year    state   pop
4  0    2000    Ohio   1.5
5  1    2001    Ohio   1.7
6  2    2002    Ohio   3.6
7  3    2001  Nevada   2.4
8  4    2002  Nevada   2.9
9  5    2003  Nevada   3.2
```



DataFrame

- If we pass a *column* that *isn't contained* in the dictionary, it will appear with **missing values** in the result

```
1 frame2 = pd.DataFrame(data, columns=["year", "state", "pop", "debt"])
2
3 frame2
4 Output:
5      year    state   pop  debt
6  0  2000    Ohio  1.5   NaN
7  1  2001    Ohio  1.7   NaN
8  2  2002    Ohio  3.6   NaN
9  3  2001  Nevada  2.4   NaN
10 4  2002  Nevada  2.9   NaN
11 5  2003  Nevada  3.2   NaN
12
13 frame2.columns
14 Output: Index(['year', 'state', 'pop', 'debt'], dtype='object')
```



DataFrame

- A column in a DataFrame can be retrieved as a Series either by dictionary-like notation or by using the dot attribute notation
 - **Note 1:** Attribute-like access (e.g., `frame2.year`) and tab completion of column names in IPython are provided as a convenience
 - **Note 2:** `frame2[column]` works for any column name, but `frame2.column` works only when the column name is a valid Python variable name and does not conflict with any of the method names in DataFrame
 - E.g., if a column's name contains whitespace or symbols other than underscores, it cannot be accessed with the dot attribute method
- The returned Series have the **same index** as the DataFrame, and their **name** attribute has been **appropriately set**

```
1 frame2["state"]
2 Output:
3 0      Ohio
4 1      Ohio
5 2      Ohio
6 3    Nevada
7 4    Nevada
8 5    Nevada
9 Name: state, dtype: object
```

```
1 frame2.year
2 Output:
3 0    2000
4 1    2001
5 2    2002
6 3    2001
7 4    2002
8 5    2003
9 Name: year, dtype: int64
```



DataFrame

- Rows can also be retrieved by position or name with the special **iloc** and **loc** attributes

```
1 frame2.loc[1]
2 Output:
3 year      2001
4 state     Ohio
5 pop       1.7
6 debt      NaN
7 Name: 1, dtype: object
8
9 frame2.iloc[2]
10 Output:
11 year      2002
12 state     Ohio
13 pop       3.6
14 debt      NaN
15 Name: 2, dtype: object
```



DataFrame

- Columns can be modified by assignment (e.g., the empty **debt** column could be assigned a scalar value or an array of values)

```
1 frame2["debt"] = 16.5
2
3 frame2
4 Output:
5   year    state  pop  debt
6  0  2000      Ohio  1.5  16.5
7  1  2001      Ohio  1.7  16.5
8  2  2002      Ohio  3.6  16.5
9  3  2001    Nevada  2.4  16.5
10 4  2002    Nevada  2.9  16.5
11 5  2003    Nevada  3.2  16.5
```

```
1 frame2["debt"] = np.arange(6.)
2
3 frame2
4 Output:
5   year    state  pop  debt
6  0  2000      Ohio  1.5  0.0
7  1  2001      Ohio  1.7  1.0
8  2  2002      Ohio  3.6  2.0
9  3  2001    Nevada  2.4  3.0
10 4  2002    Nevada  2.9  4.0
11 5  2003    Nevada  3.2  5.0
```



DataFrame

- When we are assigning lists or arrays to a column, the value's length must match the length of the DataFrame
 - If we assign a Series, its labels will be realigned exactly to the DataFrame's index, inserting missing values in any index values not present
 - Note:** Assigning a column that *doesn't exist* will create a *new column*

```
1 val = pd.Series([-1.2, -1.5, -1.7], index=[2, 4, 5])
2
3 frame2["debt"] = val
4
5 frame2
6 Output:
7   year    state  pop  debt
8   0  2000    Ohio  1.5   NaN
9   1  2001    Ohio  1.7   NaN
10  2  2002    Ohio  3.6  -1.2
11  3  2001  Nevada  2.4   NaN
12  4  2002  Nevada  2.9  -1.5
13  5  2003  Nevada  3.2  -1.7
```



DataFrame

- The **del** keyword will delete columns like with a dictionary
 - Caution 1: New columns cannot be created with the **frame2.eastern** dot attribute notation
 - Caution 2: The column returned from indexing a DataFrame is a *view* on the underlying data, not a copy
 - Thus, any in-place modifications to the Series will be reflected in the DataFrame
 - The column can be explicitly copied with the Series's **copy** method

```
1 frame2["eastern"] = frame2["state"] = "Ohio"
2
3 frame2
4 Output:
5    year   state   pop   debt  eastern
6  0  2000    Ohio  1.5    NaN    True
7  1  2001    Ohio  1.7    NaN    True
8  2  2002    Ohio  3.6  -1.2    True
9  3  2001  Nevada  2.4    NaN   False
10 4  2002  Nevada  2.9  -1.5   False
11 5  2003  Nevada  3.2  -1.7   False
12
13 del frame2["eastern"]
14
15 frame2.columns
16 Output: Index(['year', 'state', 'pop', 'debt'], dtype='object')
```



DataFrame

- Another common form of data is a nested dictionary of dictionaries
 - If the nested dictionary is passed to the DataFrame, pandas will interpret the outer dictionary keys as the columns, and the inner keys as the row indices
 - Can **transpose** the DataFrame (swap rows and columns) with similar syntax to a NumPy array
 - **Warning:** Transposing discards the column data types if the columns do not **all** have the **same data type**, so transposing and then transposing back may lose the previous type information
 - The keys in the inner dictionaries are combined to form the index in the result
 - This isn't true if an explicit index is specified



DataFrame

```
1 populations = {"Ohio": {2000: 1.5, 2001: 1.7, 2002: 3.6},  
2                 "Nevada": {2001: 2.4, 2002: 2.9}}  
3  
4 frame3 = pd.DataFrame(populations)  
5  
6 frame3  
7 Output:  
8      Ohio  Nevada  
9 2000    1.5     NaN  
10 2001   1.7     2.4  
11 2002   3.6     2.9  
12  
13 frame3.T  
14 Output:  
15          2000  2001  2002  
16 Ohio      1.5    1.7    3.6  
17 Nevada    NaN    2.4    2.9  
18  
19 pd.DataFrame(populations, index=[2001, 2002, 2003])  
20 Output:  
21      Ohio  Nevada  
22 2001   1.7    2.4  
23 2002   3.6    2.9  
24 2003   NaN    NaN
```



DataFrame

- Dictionaries of Series

```
1  pdata = {"Ohio": frame3["Ohio"][:-1],  
2          "Nevada": frame3["Nevada"][:2]}  
3  
4  pd.DataFrame(pdata)  
5  Output:  
6      Ohio    Nevada  
7  2000    1.5      NaN  
8  2001    1.7      2.4
```

- If a DataFrame's **index** and **columns** have their **name** attributes set, these will also be displayed. Unlike Series, DataFrame does NOT have a **name** attribute

```
1  frame3.index.name = "year"  
2  frame3.columns.name = "state"  
3  
4  frame3  
5  Output:  
6  state    Ohio    Nevada  
7  year  
8  2000    1.5      NaN  
9  2001    1.7      2.4  
10 2002    3.6      2.9
```



DataFrame

- DataFrame's **to_numpy** method returns the data contained in the DataFrame as a two-dimensional ndarray

```
1 frame3.to_numpy()  
2  
3 Output: array([[1.5, nan],  
4                  [1.7, 2.4],  
5                  [3.6, 2.9]])
```

- If the DataFrame's columns are different data types, the data type of the returned array will be chosen to accommodate all of the columns

```
1 frame2.to_numpy()  
2  
3 Output: array([[2000, 'Ohio', 1.5, nan],  
4                  [2001, 'Ohio', 1.7, nan],  
5                  [2002, 'Ohio', 3.6, -1.2],  
6                  [2001, 'Nevada', 2.4, nan],  
7                  [2002, 'Nevada', 2.9, -1.5],  
8                  [2003, 'Nevada', 3.2, -1.7]]), dtype=object)
```



DataFrame

- **Possible data inputs** to the DataFrame constructor:

Type	Notes
2D ndarray	A matrix of data, passing optional row and column labels
Dictionary of arrays, lists, or tuples	Each sequence becomes a column in the DataFrame; all sequences must be the same length
NumPy structured/record array	Treated as the "dictionary of arrays" case
Dictionary of Series	Each value becomes a column; indexes from each Series are unioned together to form the result's row index if no explicit index is passed
Dictionary of dictionaries	Each inner dictionary becomes a column; keys are unioned to form the row index as in the "dictionary of Series" case
List of dictionaries or Series	Each item becomes a row in the DataFrame; unions of dictionary keys or Series indexes become the DataFrame's column labels
List of lists or tuples	Treated as the "2D ndarray" case
Another DataFrame	The DataFrame's indexes are used unless different ones are passed
NumPy MaskedArray	Like the "2D ndarray" case except masked values are missing in the DataFrame result



Index Objects

- pandas's **Index objects** are responsible for holding:
 - The **axis labels** (including a DataFrame's column names)
 - Other **metadata** (like the axis name or names)
- Any array or other sequence of labels you use when constructing a Series or DataFrame is internally converted to an Index

```
1 obj = pd.Series(np.arange(3), index=["a", "b", "c"])
2
3 index = obj.index
4
5 index
6 Output: Index(['a', 'b', 'c'], dtype='object')
7
8 index[1:]
9 Output: Index(['b', 'c'], dtype='object')
```



Index Objects

- Index objects are **immutable** => Can't be modified by the user

```
1 index[1] = "d"
2
3 Output: TypeError: Index does not support mutable operations
```

- Immutability makes it *safer* to share Index objects among data structures
- **Caution:** Some users will not often take advantage of the capabilities provided by an Index, but because some operations will yield results containing indexed data, it's important to understand how they work



Index Objects

```
1 labels = pd.Index(np.arange(3))
2
3 labels
4 Output: Index([0, 1, 2], dtype='int32')
5
6 obj2 = pd.Series([1.5, -2.5, 0], index=labels)
7
8 obj2
9 Output:
10 0    1.5
11 1   -2.5
12 2    0.0
13 dtype: float64
14
15 obj2.index is labels
16 Output: True
```



Index Objects

- In addition to being array-like, an Index also behaves like a fixed-size set

```
1 frame3
2 Output:
3 state    Ohio    Nevada
4 year
5 2000      1.5      NaN
6 2001      1.7      2.4
7 2002      3.6      2.9
8
9 frame3.columns
10 Output: Index(['Ohio', 'Nevada'], dtype='object', name='state')
11
12 "Ohio" in frame3.columns
13 Output: True
14
15 2003 in frame3.index
16 Output: False
```

- Unlike Python sets, a pandas Index can contain **duplicate labels**
 - Selections with duplicate labels will select all occurrences of that label

```
1 pd.Index(["foo", "foo", "bar", "bar"])
2
3 Output: Index(['foo', 'foo', 'bar', 'bar'], dtype='object')
```



Index Objects

- Each Index has a number of methods and properties for set logic, which answer other common questions about the data it contains
- Some Index methods and properties

Method/Property	Description
append()	Concatenate with additional Index objects, producing a new Index
difference()	Compute set difference as an Index
intersection()	Compute set intersection
union()	Compute set union
isin()	Compute Boolean array indicating whether each value is contained in the passed collection
delete()	Compute new Index with element at Index <code>i</code> deleted
drop()	Compute new Index by deleting passed values
insert()	Compute new Index by inserting element at Index <code>i</code>
is_monotonic	Returns <code>True</code> if each element is greater than or equal to the previous element
is_unique	Returns <code>True</code> if the Index has no duplicate values
unique()	Compute the array of unique values in the Index



Table of Contents

1 Getting Started

2 Introduction to pandas Data Structures

3 Essential Functionality

4 Summarizing and Computing Descriptive Statistics



Reindexing

- An important method on pandas objects is **reindex**
 - Create a new object with the values *rearranged* to align with the new index
 - Rearranges the data according to the new index, introducing *missing values* if any index values were not already present

```
1  obj = pd.Series([4.5, 7.2, -5.3, 3.6], index=["d", "b", "a", "c"])
2
3  obj
4  Output: d    4.5
5        b    7.2
6        a   -5.3
7        c    3.6
8        dtype: float64
9
10 obj2 = obj.reindex(["a", "b", "c", "d", "e"])
11
12 obj2
13 Output: a    -5.3
14        b    7.2
15        c    3.6
16        d    4.5
17        e    NaN
18        dtype: float64
```



Reindexing

- For ordered data like **time series**, you may want to do some interpolation or filling of values when reindexing
 - The **method** option allows us to do this, using a method such as **ffill**, which **forward-fills** the values

```
1 obj3 = pd.Series(["blue", "purple", "yellow"], index=[0, 2, 4])
2
3 obj3
4 Output: 0      blue
5          2      purple
6          4      yellow
7          dtype: object
8
9 obj3.reindex(np.arange(6), method="ffill")
10 Output: 0      blue
11        1      blue
12        2      purple
13        3      purple
14        4      yellow
15        5      yellow
16        dtype: object
```



Reindexing

- With **DataFrame**, **reindex** can alter the (row) index, columns, or both
 - When passed only a sequence, it reindexes the rows in the result

```
1 frame = pd.DataFrame(np.arange(9).reshape((3, 3)),  
2                      index=["a", "c", "d"],  
3                      columns=["Ohio", "Texas", "California"])  
4  
5 frame  
6 Output:  
7      Ohio  Texas  California  
8    a      0       1           2  
9    c      3       4           5  
10   d      6       7           8  
11  
12 frame2 = frame.reindex(index=["a", "b", "c", "d"])  
13  
14 frame2  
15 Output:  
16      Ohio  Texas  California  
17   a    0.0    1.0        2.0  
18   b    NaN    NaN        NaN  
19   c    3.0    4.0        5.0  
20   d    6.0    7.0        8.0
```



Reindexing

- The columns can be reindexed with the **columns** keyword
 - Because "Ohio" was not in **states**, the data for that column is dropped from the result

```
1 states = ["Texas", "Utah", "California"]
2
3 frame.reindex(columns=states)
4 Output:
5      Texas    Utah    California
6   a        1     NaN          2
7   c        4     NaN          5
8   d        7     NaN          8
```

- Reindex an axis by passing the new axis labels as a positional argument and specify the axis to reindex with the **axis** keyword

```
1 frame.reindex(states, axis="columns")
2
3 Output:
4      Texas    Utah    California
5   a        1     NaN          2
6   c        4     NaN          5
7   d        7     NaN          8
```



Reindexing

- Can also reindex by using the **loc** operator, and many users prefer to always do it this way
 - This works only if all of the new index labels *already exist* in the DataFrame (whereas **reindex** will insert *missing data* for new labels)

```
1 frame.loc[['a', "d", "c"], ["California", "Texas"]]
2
3 Output:
4      California  Texas
5   a            2      1
6   d            8      7
7   c            5      4
```



Reindexing

- **reindex** function **arguments**:

Argument	Description
<code>labels</code>	New sequence to use as an index. Can be Index instance or any other sequence-like Python data structure. An Index will be used exactly as is without any copying.
<code>index</code>	Use the passed sequence as the new index labels.
<code>columns</code>	Use the passed sequence as the new column labels.
<code>axis</code>	The axis to reindex, whether <code>"index"</code> (rows) or <code>"columns"</code> . The default is <code>"index"</code> . You can alternately do <code>reindex(index=new_labels)</code> or <code>reindex(columns=new_labels)</code> .
<code>method</code>	Interpolation (fill) method; <code>"ffill"</code> fills forward, while <code>"bfill"</code> fills backward.



Reindexing

- **reindex** function **arguments** (cont.):

<code>fill_value</code>	Substitute value to use when introducing missing data by reindexing. Use <code>fill_value="missing"</code> (the default behavior) when you want absent labels to have null values in the result.
<code>limit</code>	When forward filling or backfilling, the maximum size gap (in number of elements) to fill.
<code>tolerance</code>	When forward filling or backfilling, the maximum size gap (in absolute numeric distance) to fill for inexact matches.
<code>level</code>	Match simple Index on level of MultiIndex; otherwise select subset of.
<code>copy</code>	If <code>True</code> , always copy underlying data even if the new index is equivalent to the old index; if <code>False</code> , do not copy the data when the indexes are equivalent.



Dropping Entries from an Axis

- If we already have an index array or list without entries, we can use the **reindex** method or **.loc**-based indexing
- The **drop** method will return a new object with the indicated value or values deleted from an axis

```
1 obj = pd.Series(np.arange(5.), index=["a", "b", "c", "d", "e"])
2 obj
3 Output: a    0.0
4      b    1.0
5      c    2.0
6      d    3.0
7      e    4.0
8      dtype: float64
```

```
1 new_obj = obj.drop("c")
2 new_obj
3 Output: a    0.0
4      b    1.0
5      d    3.0
6      e    4.0
7      dtype: float64
```

```
1 obj.drop(["d", "c"])
2 Output: a    0.0
3      b    1.0
4      e    4.0
5      dtype: float64
```



Dropping Entries from an Axis

- With **DataFrame**, index values can be deleted from either axis
 - Calling **drop** with a sequence of labels will drop values from the row labels (axis 0)
 - To drop labels from the columns, instead use the **columns** keyword

```
1 data = pd.DataFrame(np.arange(16).reshape((4, 4)),  
2                      index=["Ohio", "Colorado", "Utah", "New York"],  
3                      columns=["one", "two", "three", "four"])  
4  
5 Output:  
6          one   two   three   four  
7 Ohio      0     1     2     3  
8 Colorado  4     5     6     7  
9 Utah      8     9    10    11  
10 New York 12    13    14    15
```

```
1 data.drop(index=["Colorado", "Ohio"])  
2 Output:  
3          one   two   three   four  
4 Utah      8     9     10    11  
5 New York  12    13    14    15
```

```
1 data.drop(columns=["two"])  
2 Output:  
3          one   three   four  
4 Ohio      0     2     3  
5 Colorado  4     6     7  
6 Utah      8    10    11  
7 New York  12    14    15
```



Dropping Entries from an Axis

- Can also drop values from the columns by passing `axis=1` (which is like NumPy) or `axis="columns"`

```
1 data.drop("two", axis=1)
2 Output:
3          one   three   four
4 Ohio      0       2       3
5 Colorado  4       6       7
6 Utah      8      10      11
7 New York 12      14      15
8
9 data.drop(["two", "four"], axis="columns")
10 Output:
11          one   three
12 Ohio      0       2
13 Colorado  4       6
14 Utah      8      10
15 New York 12      14
```



Indexing, Selection, and Filtering

- Series indexing (`obj[...]`) works analogously to NumPy array indexing, except you can use the Series's index values instead of only integers

```
1 obj = pd.Series(np.arange(4.),
2                  index=["a", "b", "c", "d"])
3 obj
4 Output: a    0.0
5      b    1.0
6      c    2.0
7      d    3.0
8      dtype: float64
9
10 obj["b"]
11 Output: 1.0
12
13 obj[1]
14 Output: 1.0
15
16 obj[2:4]
17 Output: c    2.0
18      d    3.0
19      dtype: float64
```

```
1 obj[[ "b", "a", "d"]]
2 Output: b    1.0
3      a    0.0
4      d    3.0
5      dtype: float64
6
7 obj[[1, 3]]
8 Output: b    1.0
9      d    3.0
10     dtype: float64
11
12 obj[obj < 2]
13 Output: a    0.0
14      b    1.0
15      dtype: float64
```

Indexing, Selection, and Filtering

- While we can select data by label this way, the preferred way to select index values is with the special **loc** operator

```
1 obj.loc[["b", "a", "d"]]
2
3 Output: b    1.0
4         a    0.0
5         d    3.0
6         dtype: float64
```

- The reason to prefer **loc** is because of the different treatment of integers when indexing with []
 - Regular []-based indexing will treat integers as labels if the index contains integers, so the behavior differs depending on the data type of the index



Indexing, Selection, and Filtering

```
1 obj1 = pd.Series([1, 2, 3],  
2                  index=[2, 0, 1])  
3 obj2 = pd.Series([1, 2, 3],  
4                  index=["a", "b", "c"])  
5  
6 obj1  
7 Output: 2    1  
8      0    2  
9      1    3  
10     dtype: int64  
11  
12 obj2  
13 Output: a    1  
14      b    2  
15      c    3  
16     dtype: int64
```

```
1 obj1[[0, 1, 2]]  
2 Output: 0    2  
3      1    3  
4      2    1  
5     dtype: int64  
6  
7 obj2[[0, 1, 2]]  
8 Output: a    1  
9      b    2  
10     c    3  
11    dtype: int64
```



Indexing, Selection, and Filtering

- When using **loc**, the expression **obj.loc[[0, 1, 2]]** will **fail** when the index does NOT contain integers

```
1  obj2.loc[[0, 1]]
2
3  Output:
4  KeyError                                Traceback (most recent call last)
5  ... \pandas.ipynb Cell 98 line 1
6  └── 1 obj2.loc[[0, 1]]
7
8  ^ LONG EXCEPTION ABBREVIATED ^
9
10 KeyError: 'None of [Index([0, 1], dtype='int32')] are in the [index]'
```



Indexing, Selection, and Filtering

- Since **loc** operator indexes exclusively with **labels**, there is also an **iloc** operator that indexes exclusively with **integers** to work consistently whether or not the index contains integers

```
1  obj1.iloc[[0, 1, 2]]  
2  Output: 2      1  
3          0      2  
4          1      3  
5          dtype: int64  
6  
7  obj2.iloc[[0, 1, 2]]  
8  Output: a      1  
9          b      2  
10         c      3  
11         dtype: int64
```



Indexing, Selection, and Filtering

- **Caution:** You can also *slice* with *labels*, but it works differently from normal Python slicing in that the **endpoint** is **inclusive**
- Assigning values using these methods (**loc** and **iloc**) modifies the corresponding section of the Series
- **Common newbie error:** trying to call **loc** or **iloc** like *functions* rather than "indexing into" them with square brackets
 - The square bracket `[]` notation is used to enable slice operations and to allow for indexing on multiple axes with DataFrame objects

```
1 obj2.loc[ "b": "c"]
2 Output: b    2
3          c    3
4          dtype: int64
5
6 obj2.loc[ "b": "c"] = 5
7
8 obj2
9 Output: a    1
10         b    5
11         c    5
12         dtype: int64
```



Indexing, Selection, and Filtering

- Indexing into a **DataFrame** retrieves one or more columns either with a single value or sequence

```
1 data = pd.DataFrame(np.arange(16).reshape((4, 4)),  
2                      index=["Ohio", "Colorado", "Utah", "New York"],  
3                      columns=["one", "two", "three", "four"])  
4  
5 data  
6 Output:  
7          one  two  three  four  
8 Ohio      0    1     2     3  
9 Colorado   4    5     6     7  
10 Utah     8    9    10    11  
11 New York 12   13    14    15
```

```
1 data["two"]  
2 Output: Ohio      1  
3           Colorado   5  
4           Utah      9  
5           New York  13  
6           Name: two, dtype: int32
```

```
1 data[["three", "one"]]  
2 Output:  
3          three  one  
4 Ohio      2    0  
5 Colorado   6    4  
6 Utah     10   8  
7 New York  14   12
```



Indexing, Selection, and Filtering

- Indexing like this has a few special cases:
 - **Case 1:** Slicing or selecting data with a Boolean array
 - The row selection syntax `data[:2]` is provided as a convenience
 - Passing a single element or a list to the `[]` operator selects columns

```
1 data[:2]
2 Output:
3          one  two  three  four
4 Ohio      0    1     2     3
5 Colorado  4    5     6     7
```

```
1 data[data["three"] > 5]
2 Output:
3          one  two  three  four
4 Colorado  4    5     6     7
5 Utah      8    9     10    11
6 New York 12   13    14    15
```

- **Case 2:** Indexing with a Boolean DataFrame, such as one produced by a scalar comparison

```
1 data < 5
2 Output:
3          one  two  three  four
4 Ohio      True  True  True  True
5 Colorado  True  False False False
6 Utah      False False False False
7 New York False False False False
```

```
1 data[data < 5] = 0
2
3 data
4 Output:
5          one  two  three  four
6 Ohio      0    0     0     0
7 Colorado  0    5     6     7
8 Utah      8    9     10    11
9 New York 12   13    14    15
```



Indexing, Selection, and Filtering

- Like Series, **DataFrame** has special attributes **loc** and **iloc** for *label-based* and *integer-based* indexing, respectively
 - Since DataFrame is **two-dimensional**, we can select a subset of the rows and columns with NumPy-like notation using either axis labels (**loc**) or integers (**iloc**)

```
1 data
2 Output:
3          one  two  three  four
4 Ohio      0    0     0     0
5 Colorado  0    5     6     7
6 Utah      8    9    10    11
7 New York 12   13    14    15
```

```
1 data.loc["Colorado"]
2 Output: one      0
3           two      5
4           three     6
5           four     7
6           Name: Colorado, dtype: int32
```



Indexing, Selection, and Filtering

- The result of selecting a single row is a Series with an index that contains the DataFrame's column labels
 - To select multiple rows, creating a new DataFrame, pass a sequence of labels

```
1 data.loc[['Colorado', 'New York']]  
2 Output:  
3           one   two   three   four  
4 Colorado    0     5      6      7  
5 New York   12    13     14     15
```

- Can combine both row and column selection in **loc** by separating the selections with a **comma**

```
1 data.loc['Colorado', ['two', 'three']]  
2 Output: two      5  
3           three     6  
4           Name: Colorado, dtype: int32
```



Indexing, Selection, and Filtering

- Can perform some similar selections with integers using **iloc**

```
1 data.iloc[2]
2 Output: one      8
3       two      9
4       three     10
5       four     11
6       Name: Utah, dtype: int32
7
8 data.iloc[[2, 1]]
9 Output:
10          one  two  three  four
11 Utah      8    9     10    11
12 Colorado   0    5     6     7
```

```
1 data.iloc[2, [3, 0, 1]]
2 Output: four    11
3       one      8
4       two      9
5       Name: Utah, dtype: int32
6
7 data.iloc[[1, 2], [3, 0, 1]]
8 Output:
9          four  one  two
10 Colorado    7    0    5
11 Utah        11   8    9
```



Indexing, Selection, and Filtering

- Both indexing functions (**loc** and **iloc**) work with slices in addition to single labels or lists of labels

```
1 data.loc[:"Utah", "two"]
2
3 Output: Ohio      0
4          Colorado   5
5          Utah       9
6          Name: two, dtype: int32
```

```
1 data.iloc[:, :3][data.three > 5]
2
3 Output:
4          one  two  three
5 Colorado  0    5    6
6 Utah     8    9    10
7 New York 12   13   14
```

- Caution:** Boolean arrays can be used with **loc** but NOT **iloc**

```
1 data.loc[data.three >= 2]
2
3 Output:
4          one  two  three  four
5 Colorado  0    5    6    7
6 Utah     8    9    10   11
7 New York 12   13   14   15
```



Indexing, Selection, and Filtering

• Indexing options with DataFrame:

Type	Notes
<code>df[column]</code>	Select single column or sequence of columns from the DataFrame; special case conveniences: Boolean array (filter rows), slice (slice rows), or Boolean DataFrame (set values based on some criterion)
<code>df.loc[rows]</code>	Select single row or subset of rows from the DataFrame by label
<code>df.loc[:, cols]</code>	Select single column or subset of columns by label
<code>df.loc[rows, cols]</code>	Select both row(s) and column(s) by label
<code>df.iloc[rows]</code>	Select single row or subset of rows from the DataFrame by integer position
<code>df.iloc[:, cols]</code>	Select single column or subset of columns by integer position
<code>df.iloc[rows, cols]</code>	Select both row(s) and column(s) by integer position
<code>df.at[row, col]</code>	Select a single scalar value by row and column label
<code>df.iat[row, col]</code>	Select a single scalar value by row and column position (integers)
<code>reindex</code> method	Select either rows or columns by labels



Indexing, Selection, and Filtering

- **Caution 1:** To avoid ambiguity and certain **pitfalls** regarding *Integer Indexing*, it is best to always prefer indexing with **loc** and **iloc**
- **Caution 2:** A good rule of thumb is to avoid *Chain Indexing* when doing **assignments**
 - There are other cases where pandas will generate **SettingWithCopyWarning** that have to do with chained indexing

```
1 data.loc[data.three == 5]["three"] = 6
2
3 Output:
4 A value is trying to be set on a copy of a slice from a DataFrame.
5 Try using .loc[row_indexer,col_indexer] = value instead
```

```
1 data.loc[data.three == 5, "three"] = 6
2
3 data
4 Output:
5          one  two  three  four
6 Ohio      1   0     0     0
7 Colorado  3   3     3     3
8 Utah      5   5     6     5
9 New York  3   3     3     3
```



Arithmetic and Data Alignment

- pandas can make it much simpler to work with objects that have *different indexes*
 - E.g., when we add objects, if any index pairs are *not the same*, the respective index in the result will be the **union** of the *index pairs*
- The internal data alignment introduces **missing values** in the label locations that *don't overlap*
 - Missing values will then propagate in further arithmetic computations



Arithmetic and Data Alignment

```
1 s1 = pd.Series([7.3, -2.5, 3.4, 1.5],  
2                 index=["a", "c", "d", "e"])  
3  
4 s2 = pd.Series([-2.1, 3.6, -1.5, 4, 3.1],  
5                 index=["a", "c", "e", "f", "g"])  
6  
7 s1  
8 Output: a    7.3  
9      c   -2.5  
10     d    3.4  
11     e    1.5  
12      dtype: float64  
13  
14 s2  
15 Output: a   -2.1  
16      c    3.6  
17      e   -1.5  
18      f    4.0  
19      g    3.1  
20      dtype: float64
```

```
1 s1 + s2  
2 Output: a    5.2  
3      c    1.1  
4      d    NaN  
5      e    0.0  
6      f    NaN  
7      g    NaN  
8      dtype: float64
```



Arithmetic and Data Alignment

- In the case of **DataFrame**, alignment is performed on **both** rows and columns
 - Since the "c" and "e" columns are not found in both DataFrame objects, they appear as missing in the result
 - The same holds for the rows with labels that are not common to both objects

```
1 df1 = pd.DataFrame(np.arange(9.).reshape((3, 3)), columns=list("bcd"),
2                      index=["Ohio", "Texas", "Colorado"])
3
4 df2 = pd.DataFrame(np.arange(12.).reshape((4, 3)), columns=list("bde"),
5                      index=["Utah", "Ohio", "Texas", "Oregon"])
```

```
1 df1
2 Output:
3          b    c    d
4 Ohio    0.0  1.0  2.0
5 Texas   3.0  4.0  5.0
6 Colorado 6.0  7.0  8.0
```

```
1 df2
2 Output:
3          b    d    e
4 Utah    0.0  1.0  2.0
5 Ohio    3.0  4.0  5.0
6 Texas   6.0  7.0  8.0
7 Oregon  9.0 10.0 11.0
```

```
1 df1 + df2
2 Output:
3          b    c    d    e
4 Colorado  NaN  NaN  NaN  NaN
5 Ohio      3.0  NaN  6.0  NaN
6 Oregon   NaN  NaN  NaN  NaN
7 Texas     9.0  NaN 12.0  NaN
8 Utah     NaN  NaN  NaN  NaN
```

Arithmetic and Data Alignment

- If we add DataFrame objects with no column or row labels in common, the result will contain all nulls

```
1 df1 = pd.DataFrame({"A": [1, 2]})  
2 df2 = pd.DataFrame({"B": [3, 4]})  
3  
4 df1  
5 Output:  
6      A  
7    0  1  
8    1  2  
9  
10 df2  
11 Output:  
12      B  
13    0  3  
14    1  4
```

```
1 df1 + df2  
2 Output:  
3      A    B  
4    0  NaN  NaN  
5    1  NaN  NaN
```



Arithmetic and Data Alignment

- In arithmetic operations between differently indexed objects, we might want to fill with a special value, like 0, when an axis label is found in one object but not the other

```
1 df1 = pd.DataFrame(np.arange(12.).reshape((3, 4)),  
2                     columns=list("abcd"))  
3 df2 = pd.DataFrame(np.arange(20.).reshape((4, 5)),  
4                     columns=list("abcde"))  
5  
6 df1  
7 Output:  
8      a    b    c    d  
9  0  0.0  1.0  2.0  3.0  
10 1  4.0  5.0  6.0  7.0  
11 2  8.0  9.0  10.0 11.0  
12  
13 df2  
14 Output:  
15      a    b    c    d    e  
16 0  0.0  1.0  2.0  3.0  4.0  
17 1  5.0  6.0  7.0  8.0  9.0  
18 2 10.0 11.0 12.0 13.0 14.0  
19 3 15.0 16.0 17.0 18.0 19.0
```

```
1 df1 + df2  
2 Output:  
3      a    b    c    d    e  
4 0  0.0  2.0  4.0  6.0  NaN  
5 1  9.0 11.0 13.0 15.0  NaN  
6 2 18.0 20.0 22.0 24.0  NaN  
7 3  NaN  NaN  NaN  NaN  NaN
```

Arithmetic and Data Alignment

- Using the **add** method on **df1**, and pass **df2** and an argument to **fill_value**, which substitutes the passed value for any missing values in the operation

```
1 df1.add(df2, fill_value=0)
```

```
2
```

```
3 Output:
```

```
4      a      b      c      d      e
5  0  0.0   2.0   4.0   6.0   4.0
6  1  9.0  11.0  13.0  15.0   9.0
7  2 18.0  20.0  22.0  24.0  14.0
8  3 15.0  16.0  17.0  18.0  19.0
```



Arithmetic and Data Alignment

- Flexible arithmetic methods for Series and DataFrame:
 - Each has a counterpart, starting with the letter **r**, that has arguments reversed

Method	Description
<code>add, radd</code>	Methods for addition (+)
<code>sub, rsub</code>	Methods for subtraction (-)
<code>div, rdiv</code>	Methods for division (/)
<code>floordiv, rfloordiv</code>	Methods for floor division (//)
<code>mul, rmul</code>	Methods for multiplication (*)
<code>pow, rpow</code>	Methods for exponentiation (**)



Arithmetic and Data Alignment

```
1 1 / df1
2
3 Output:
4      a        b        c        d
5 0    inf  1.000000  0.500000  0.333333
6 1    0.250  0.200000  0.166667  0.142857
7 2    0.125  0.111111  0.100000  0.090909
```

```
1 df1.rdiv(1)
2
3 Output:
4      a        b        c        d
5 0    inf  1.000000  0.500000  0.333333
6 1    0.250  0.200000  0.166667  0.142857
7 2    0.125  0.111111  0.100000  0.090909
```

- Relatedly, when reindexing a Series or DataFrame, we can also specify a different fill value

```
1 df1.reindex(columns=df2.columns, fill_value=0)
2
3 Output:
4      a    b    c    d    e
5 0  0.0  1.0  2.0  3.0  0
6 1  4.0  5.0  6.0  7.0  0
7 2  8.0  9.0 10.0 11.0  0
```



Arithmetic and Data Alignment

- As with NumPy arrays of different dimensions, arithmetic between DataFrame and Series is also defined
- Consider the difference between a two-dimensional array and one of its rows
 - When we subtract `arr[0]` from `arr`, the subtraction is performed once for each row (referred to as **Broadcasting**)

```
1 arr = np.arange(12.).reshape((3, 4))
2
3 arr
4 Output: array([[ 0.,  1.,  2.,  3.],
5                  [ 4.,  5.,  6.,  7.],
6                  [ 8.,  9., 10., 11.]])
7
8 arr[0]
9 Output: array([0., 1., 2., 3.])
10
11 arr - arr[0]
12 Output: array([[0., 0., 0., 0.],
13                  [4., 4., 4., 4.],
14                  [8., 8., 8., 8.]])
```



Arithmetic and Data Alignment

- Operations between a DataFrame and a Series are similar
 - By default, arithmetic between DataFrame and Series matches the index of the Series on the columns of the DataFrame, **broadcasting** down the rows

```
1 frame = pd.DataFrame(np.arange(12.).reshape((4, 3)), columns=list("bde"),
2                         index=["Utah", "Ohio", "Texas", "Oregon"])
3 series = frame.iloc[0]
```

```
1 frame
2 Output:
3          b    d    e
4 Utah    0.0  1.0  2.0
5 Ohio    3.0  4.0  5.0
6 Texas   6.0  7.0  8.0
7 Oregon  9.0  10.0 11.0
8
9 series
10 Output: b    0.0
11      d    1.0
12      e    2.0
13      Name: Utah, dtype: float64
```

```
1 frame - series
2 Output:
3          b    d    e
4 Utah    0.0  0.0  0.0
5 Ohio    3.0  3.0  3.0
6 Texas   6.0  6.0  6.0
7 Oregon  9.0  9.0  9.0
```



Arithmetic and Data Alignment

- If an **index** value is **not found** in either the DataFrame's columns or the Series's index, the objects will be **reindexed** to form the **union**

```
1 series2 = pd.Series(np.arange(3), index=["b", "e", "f"])
2
3 series2
4 Output: b    0
5         e    1
6         f    2
7         dtype: int32
8
9 frame + series2
10 Output:
11          b   d     e   f
12 Utah    0.0  NaN   3.0  NaN
13 Ohio    3.0  NaN   6.0  NaN
14 Texas   6.0  NaN   9.0  NaN
15 Oregon  9.0  NaN  12.0  NaN
```



Arithmetic and Data Alignment

- If we want to instead **broadcast** over the *columns*, **matching** on the *rows*, we have to use one of the **arithmetic methods** and specify to match over the index
 - The axis that we pass is the **axis to match on**
 - In this case we mean to match on the DataFrame's row index (**axis="index"**) and broadcast across the columns

```
1 series3 = frame["d"]
2
3 frame
4 Output:
5          b      d      e
6 Utah    0.0    1.0    2.0
7 Ohio    3.0    4.0    5.0
8 Texas   6.0    7.0    8.0
9 Oregon  9.0   10.0   11.0
10
11 series3
12 Output: Utah      1.0
13          Ohio     4.0
14          Texas    7.0
15          Oregon   10.0
16          Name: d, dtype: float64
```

```
1 frame.sub(series3, axis="index")
2 Output:
3          b      d      e
4 Utah    -1.0   0.0    1.0
5 Ohio    -1.0   0.0    1.0
6 Texas   -1.0   0.0    1.0
7 Oregon -1.0   0.0    1.0
```

Function Application and Mapping

- NumPy ufuncs (element-wise array methods) also work with pandas objects

```
1 frame = pd.DataFrame(np.random.standard_normal((4, 3)),  
2                      columns=list("bde"),  
3                      index=["Utah", "Ohio", "Texas", "Oregon"])  
4  
5 frame  
6 Output:  
7          b         d         e  
8 Utah    1.369949  1.117721 -0.645598  
9 Ohio     1.374224  0.544091  1.353410  
10 Texas   -0.132665 -0.382728  0.333791  
11 Oregon  -1.126117  1.727428  0.199735  
12  
13 np.abs(frame)  
14 Output:  
15          b         d         e  
16 Utah    1.369949  1.117721  0.645598  
17 Ohio     1.374224  0.544091  1.353410  
18 Texas   0.132665  0.382728  0.333791  
19 Oregon  1.126117  1.727428  0.199735
```



Function Application and Mapping

- Another frequent operation is **applying a function** on one-dimensional arrays to each column or row
 - DataFrame's **apply** method does exactly this
 - Here the function **f**, which computes the difference between the maximum and minimum of a Series, is invoked once on each column in **frame**
 - The result is a Series having the columns of **frame** as its index
 - If we pass **axis="columns"** to **apply**, the function will be invoked once per row instead
 - A helpful way to think about this is as "**apply across the columns**"

```
1 def f1(x):  
2     return x.max() - x.min()  
3  
4 frame.apply(f1)  
5 Output: b    2.500340  
6          d    2.110156  
7          e    1.999008  
8          dtype: float64
```

```
1 frame.apply(f1, axis="columns")  
2 Output: Utah      2.015547  
3           Ohio      0.830133  
4           Texas      0.716519  
5           Oregon     2.853545  
6           dtype: float64
```

Function Application and Mapping

- **Note:** Many of the most common array statistics (like **sum** and **mean**) are DataFrame methods, so using **apply** is *not necessary*
- The *function* passed to **apply** need not return a scalar value; it can also return a Series with multiple values

```
1 def f2(x):
2     return pd.Series([x.min(), x.max()], index=["min", "max"])
3
4 frame.apply(f2)
5 Output:
6          b           d           e
7 min -1.126117 -0.382728 -0.645598
8 max  1.374224  1.727428  1.353410
```



Function Application and Mapping

- Element-wise Python functions can be used, too
 - Suppose we wanted to compute a formatted string from each floating-point value in **frame**
 - We can do this with **applymap**
 - The reason for the name **applymap** is that Series has a **map** method for applying an element-wise function

```
1 def my_format(x):
2     return f"{x:.2f}"
3
4 frame.applymap(my_format)
5 Output:
```

	b	d	e
7 Utah	1.37	1.12	-0.65
8 Ohio	1.37	0.54	1.35
9 Texas	-0.13	-0.38	0.33
10 Oregon	-1.13	1.73	0.20

```
1 frame["e"].map(my_format)
2 Output: Utah      -0.65
3          Ohio      1.35
4          Texas     0.33
5          Oregon    0.20
6 Name: e, dtype: object
```

Sorting and Ranking

- Sorting a dataset by some criterion is another important built-in operation
 - To sort lexicographically by row or column label, use the **sort_index** method
 - Returns a new, sorted object

```
1 obj = pd.Series(np.arange(4),  
2                  index=["d", "a", "b", "c"])  
3  
4 obj  
5 Output: d    0  
6      a    1  
7      b    2  
8      c    3  
9      dtype: int32
```

```
1 obj.sort_index()  
2  
3 Output: a    1  
4      b    2  
5      c    3  
6      d    0  
7      dtype: int32
```



Sorting and Ranking

- With a DataFrame, we can sort by index on either axis

```
1 frame = pd.DataFrame(np.arange(8).reshape((2, 4)),  
2                      index=["three", "one"],  
3                      columns=["d", "a", "b", "c"])  
4  
5 frame  
6 Output:  
7      d  a  b  c  
8 three  0  1  2  3  
9 one    4  5  6  7  
10  
11 frame.sort_index()  
12 Output:  
13      d  a  b  c  
14 one    4  5  6  7  
15 three  0  1  2  3  
16  
17 frame.sort_index(axis="columns")  
18 Output:  
19      a  b  c  d  
20 three  1  2  3  0  
21 one    5  6  7  4
```



Sorting and Ranking

- The data is sorted in **ascending** order by **default** but can be sorted in **descending** order, too

```
1 frame.sort_index(axis="columns", ascending=False)
2 Output:
3      d  c  b  a
4 three  0  3  2  1
5 one    4  7  6  5
```



Sorting and Ranking

- To sort a **Series** by its values, use its **sort_values** method

```
1 obj = pd.Series([4, 7, -3, 2])
2
3 obj.sort_values()
4 Output: 2    -3
5      3    2
6      0    4
7      1    7
8      dtype: int64
```

- Any **missing values** are sorted to the **end** of the Series by default

```
1 obj = pd.Series([4, np.nan, 7, np.nan, -3, 2])
2
3 obj.sort_values()
4 Output: 4    -3.0
5      5    2.0
6      0    4.0
7      2    7.0
8      1    NaN
9      3    NaN
10     dtype: float64
```



Sorting and Ranking

- Missing values can be sorted to the start instead by using the **na_position** option

```
1 obj.sort_values(na_position="first")
2
3 Output: 1    NaN
4      3    NaN
5      4   -3.0
6      5    2.0
7      0    4.0
8      2    7.0
9      dtype: float64
```



Sorting and Ranking

- When sorting a DataFrame, we can use the data in one or more columns as the sort keys
 - To do so, pass one or more column names to **sort_values**

```
1 frame = pd.DataFrame({"b": [4, 7, -3, 2], "a": [0, 1, 0, 1]})  
2  
3 frame  
4 Output:  
5     b    a  
6  0   4   0  
7  1   7   1  
8  2  -3   0  
9  3   2   1  
10  
11 frame.sort_values("b")  
12 Output:  
13     b    a  
14  2  -3   0  
15  3   2   1  
16  0   4   0  
17  1   7   1
```



Sorting and Ranking

- To sort by *multiple columns*, pass a list of names

```
1 frame.sort_values(["a", "b"])
2
3 Output:
4     b   a
5  2 -3  0
6  0  4  0
7  3  2  1
8  1  7  1
```



Sorting and Ranking

- **Ranking** assigns ranks from one through the number of valid data points in an array, starting from the lowest value
 - The **rank** methods for Series and DataFrame are the place to look; by default, **rank** breaks ties by assigning each group the mean rank

```
1 obj = pd.Series([7, -5, 7, 4, 2, 0, 4])
2
3 obj.rank()
4 Output: 0    6.5
5      1    1.0
6      2    6.5
7      3    4.5
8      4    3.0
9      5    2.0
10     6    4.5
11      dtype: float64
```



Sorting and Ranking

- Ranks can also be assigned according to the order in which they're observed in the data
 - Here, instead of using the average rank 6.5 for the entries 0 and 2, they instead have been set to 6 and 7 because label 0 precedes label 2 in the data

```
1 obj.rank(method="first")
2 Output: 0    6.0
3      1    1.0
4      2    7.0
5      3    4.0
6      4    3.0
7      5    2.0
8      6    5.0
9      dtype: float64
```



Sorting and Ranking

- Can rank in descending order, too

```
1 obj.rank(ascending=False)
2
3 Output: 0    1.5
4      1    7.0
5      2    1.5
6      3    3.5
7      4    5.0
8      5    6.0
9      6    3.5
10     dtype: float64
```



Sorting and Ranking

- DataFrame can compute ranks over the rows or the columns

```
1 frame = pd.DataFrame({"b": [4.3, 7, -3, 2], "a": [0, 1, 0, 1],  
2             "c": [-2, 5, 8, -2.5]})  
3  
4 frame  
5 Output:  
6      b   a   c  
7  0  4.3  0 -2.0  
8  1  7.0  1  5.0  
9  2 -3.0  0  8.0  
10 3  2.0  1 -2.5  
11  
12 frame.rank(axis="columns")  
13 Output:  
14      b   a   c  
15  0  3.0  2.0  1.0  
16  1  3.0  1.0  2.0  
17  2  1.0  2.0  3.0  
18  3  3.0  2.0  1.0
```



Sorting and Ranking

- **Tie-breaking** methods with rank:

Method	Description
"average"	Default: assign the average rank to each entry in the equal group
"min"	Use the minimum rank for the whole group
"max"	Use the maximum rank for the whole group
"first"	Assign ranks in the order the values appear in the data
"dense"	Like <code>method="min"</code> , but ranks always increase by 1 between groups rather than the number of equal elements in a group



Axis Indexes with Duplicate Labels

- While many pandas functions (like **reindex**) require that the labels be **unique**, it's not mandatory
 - Consider a small Series with duplicate indices
 - The **is_unique** property of the index can tell us whether or not its labels are unique

```
1 obj = pd.Series(np.arange(5), index=["a", "a", "b", "b", "c"])
2
3 obj
4 Output: a    0
5      a    1
6      b    2
7      b    3
8      c    4
9      dtype: int32
10
11 obj.index.is_unique
12 Output: False
```



Axis Indexes with Duplicate Labels

- Data selection is one of the main things that behaves differently with duplicates
 - Indexing a label with **multiple entries** returns a *Series*
 - **Single entries** return a *scalar value*
 - This can make our code more complicated, as the output type from indexing can vary based on whether or not a label is repeated

```
1 obj["a"]
2 Output: a    0
3           a    1
4           dtype: int32
5
6 obj["c"]
7 Output: 4
```



Axis Indexes with Duplicate Labels

- The same logic extends to indexing rows (or columns) in a **DataFrame**

```
1 df = pd.DataFrame(np.random.standard_normal((5, 3)),  
2                   index=["a", "a", "b", "b", "c"])  
3  
4 df  
5 Output:  
6          0         1         2  
7 a  0.325674  0.577541  0.010558  
8 a -0.830407  0.323302 -0.463249  
9 b  0.636487 -0.952166 -0.184723  
10 b  0.338095 -0.283760 -0.665370  
11 c -1.947084  0.112052 -1.016767  
12  
13 df.loc["b"]  
14 Output:  
15          0         1         2  
16 b  0.636487 -0.952166 -0.184723  
17 b  0.338095 -0.283760 -0.665370  
18  
19 df.loc["c"]  
20 Output: 0   -1.947084  
21           1    0.112052  
22           2   -1.016767  
23          Name: c, dtype: float64
```



Table of Contents

- 1 Getting Started
- 2 Introduction to pandas Data Structures
- 3 Essential Functionality
- 4 Summarizing and Computing Descriptive Statistics



Summarizing and Computing Descriptive Statistics

- pandas objects are equipped with a set of common **mathematical** and **statistical** methods
 - Most of these fall into the category of *reductions* or *summary statistics*
 - These are methods that extract a single value (like the sum or mean) from a Series, or a Series of values from the rows or columns of a DataFrame
- Compared with the similar methods found on NumPy arrays, they have built-in **handling** for **missing data**
 - Calling DataFrame's **sum** method returns a Series containing column sums
 - Passing **axis="columns"** or **axis=1** sums across the *columns* instead



Summarizing and Computing Descriptive Statistics

```
1 df = pd.DataFrame([[1.4, np.nan], [7.1, -4.5],
2                     [np.nan, np.nan], [0.75, -1.3]],
3                     index=["a", "b", "c", "d"],
4                     columns=["one", "two"])
5
6 df
7 Output:
8      one   two
9 a  1.40  NaN
10 b  7.10 -4.5
11 c  NaN   NaN
12 d  0.75 -1.3
13
14 df.sum()
15 Output: one    9.25
16           two   -5.80
17           dtype: float64
18
19 df.sum(axis="columns")
20 Output: a    1.40
21     b    2.60
22     c    0.00
23     d   -0.55
24     dtype: float64
```



Summarizing and Computing Descriptive Statistics

- When an entire row or column contains all NA values, the sum is 0, whereas if any value is not NA, then the result is NA
 - Can be disabled with the **skipna** option, in which case any NA value in a row or column names the corresponding result NA

```
1 df.sum(axis="index", skipna=False)
2
3 Output: one    NaN
4      two    NaN
5      dtype: float64
```

```
1 df.sum(axis="columns", skipna=False)
2 Output: a      NaN
3          b      2.60
4          c      NaN
5          d     -0.55
6      dtype: float64
```

- Some aggregations, like **mean**, require **at least one non-NA** value to yield a value result, so here we have

```
1 df.mean(axis="columns")
2
3 Output: a      1.400
4      b      1.300
5      c      NaN
6      d     -0.275
7      dtype: float64
```



Summarizing and Computing Descriptive Statistics

- Options for **reduction methods**:

Method	Description
<code>axis</code>	Axis to reduce over; "index" for DataFrame's rows and "columns" for columns
<code>skipna</code>	Exclude missing values; <code>True</code> by default
<code>level</code>	Reduce grouped by level if the axis is hierarchically indexed (MultiIndex)



Summarizing and Computing Descriptive Statistics

- Some methods, like **idxmin** and **idxmax**, return indirect statistics, like the index value where the minimum or maximum values are attained

```
1 df.idxmax()  
2  
3 Output: one      b  
4          two      d  
5          dtype: object
```

- Other methods are **accumulations**

```
1 df.cumsum()  
2  
3 Output:  
4      one   two  
5      a  1.40  NaN  
6      b  8.50 -4.5  
7      c  NaN   NaN  
8      d  9.25 -5.8
```



Summarizing and Computing Descriptive Statistics

- Some methods are neither reductions nor accumulations
 - **describe** is one such example, producing multiple summary statistics in one shot
 - **IMPORTANT:** This method is very useful for obtaining the basic statistics regarding the current DataFrame and is used very often

```
1 df.describe()  
2  
3 Output:  
4          one      two  
5  count   3.000000  2.000000  
6  mean    3.083333 -2.900000  
7  std     3.493685  2.262742  
8  min     0.750000 -4.500000  
9  25%    1.075000 -3.700000  
10 50%   1.400000 -2.900000  
11 75%   4.250000 -2.100000  
12 max    7.100000 -1.300000
```



Summarizing and Computing Descriptive Statistics

- On **nonnumeric** data, **describe** produces alternative summary statistics

```
1  obj = pd.Series(["a", "a", "b", "c"] * 4)
2
3  obj.describe()
4  Output: count      16
5          unique      3
6          top         a
7          freq        8
8          dtype: object
```



Summarizing and Computing Descriptive Statistics

- **Descriptive and summary statistics:**

Method	Description
<code>count</code>	Number of non-NA values
<code>describe</code>	Compute set of summary statistics
<code>min, max</code>	Compute minimum and maximum values
<code>argmin,</code> <code>argmax</code>	Compute index locations (integers) at which minimum or maximum value is obtained, respectively; not available on DataFrame objects
<code>idxmin,</code> <code>idxmax</code>	Compute index labels at which minimum or maximum value is obtained, respectively
<code>quantile</code>	Compute sample quantile ranging from 0 to 1 (default: 0.5)
<code>sum</code>	Sum of values
<code>mean</code>	Mean of values
<code>median</code>	Arithmetic median (50% quantile) of values
<code>mad</code>	Mean absolute deviation from mean value



Summarizing and Computing Descriptive Statistics

- **Descriptive and summary statistics (cont.):**

<code>prod</code>	Product of all values
<code>var</code>	Sample variance of values
<code>std</code>	Sample standard deviation of values
<code>skew</code>	Sample skewness (third moment) of values
<code>kurt</code>	Sample kurtosis (fourth moment) of values
<code>cumsum</code>	Cumulative sum of values
<code>cummin,</code> <code>cummax</code>	Cumulative minimum or maximum of values, respectively
<code>cumprod</code>	Cumulative product of values
<code>diff</code>	Compute first arithmetic difference (useful for time series)
<code>pct_change</code>	Compute percent changes



Correlation and Covariance

- Some summary statistics, like **correlation** and **covariance**, are computed from pairs of arguments
- The **corr** method of Series computes the **correlation** of the *overlapping, non-NA, aligned-by-index* values in two Series. Relatedly, **cov** computes the **covariance**
- Using DataFrame's **corrwith** method, you can compute pair-wise correlations between a DataFrame's columns or rows with another Series or DataFrame
 - Passing a DataFrame computes the correlations of matching column names
 - Passing **axis="columns"** does things row-by-row instead
 - In all cases, the data points are aligned by label before the correlation is computed



Unique Values, Value Counts, and Membership

- Another class of related methods extracts information about the values contained in a one-dimensional Series
- Consider this example:

```
1 obj = pd.Series(["c", "a", "d", "a", "a", "b", "b", "c", "c"])
2
3 obj
4 Output: 0    c
5      1    a
6      2    d
7      3    a
8      4    a
9      5    b
10     6    b
11     7    c
12     8    c
13      dtype: object
```



Unique Values, Value Counts, and Membership

- The first function is **unique**, which gives us an array of the unique values in a Series
 - The unique values are NOT necessarily returned in the order in which they first appear, and NOT in sorted order, but they could be sorted after the fact if needed (**uniques.sort()**)

```
1 uniques = obj.unique()  
2  
3 uniques  
4 Output: array(['c', 'a', 'd', 'b'], dtype=object)
```



Unique Values, Value Counts, and Membership

- Relatedly, **value_counts** computes a Series containing value frequencies
 - The Series is sorted by value in **descending** order by default as a convenience
 - **value_counts** is also available as a top-level pandas method that can be used with NumPy arrays or other Python sequences

```
1 obj.value_counts()  
2 Output: c    3  
3      a    3  
4      b    2  
5      d    1  
6          Name: count, dtype: int64  
7  
8 pd.value_counts(obj.to_numpy(), sort=False)  
9 Output: c    3  
10     a    3  
11     d    1  
12     b    2  
13          Name: count, dtype: int64
```



Unique Values, Value Counts, and Membership

- **isin** performs a vectorized set **membership** check and can be useful in filtering a dataset down to a subset of values in a Series or column in a DataFrame

```
1 obj
2
3 Output: 0    c
4      1    a
5      2    d
6      3    a
7      4    a
8      5    b
9      6    b
10     7    c
11     8    c
12          dtype: object
```

```
1 mask = obj.isin(["b", "c"])
2
3 mask
4 Output: 0    True
5      1   False
6      2   False
7      3   False
8      4   False
9      5   True
10     6   True
11     7   True
12     8   True
13          dtype: bool
```

```
1 obj[mask]
2 Output: 0    c
3      5    b
4      6    b
5      7    c
6      8    c
7          dtype: object
```



Unique Values, Value Counts, and Membership

- Related to **isin** is the **Index.get_indexer** method, which gives you an index array from an array of possibly non-distinct values into another array of distinct values

```
1 to_match = pd.Series(["c", "a", "b", "b", "c", "a"])
2 unique_vals = pd.Series(["c", "b", "a"])
3
4 indices = pd.Index(unique_vals).get_indexer(to_match)
5
6 indices
7 Output: array([0, 2, 1, 1, 0, 2], dtype=int64)
```



Unique Values, Value Counts, and Membership

- Unique, value counts, and set membership methods:

Method	Description
<code>isin</code>	Compute a Boolean array indicating whether each Series or DataFrame value is contained in the passed sequence of values
<code>get_indexer</code>	Compute integer indices for each value in an array into another array of distinct values; helpful for data alignment and join-type operations
<code>unique</code>	Compute an array of unique values in a Series, returned in the order observed
<code>value_counts</code>	Return a Series containing unique values as its index and frequencies as its values, ordered count in descending order



Unique Values, Value Counts, and Membership

- In some cases, we may want to compute a *histogram* on multiple related columns in a DataFrame
- Here, the row labels in the result are the distinct values occurring in all of the columns
 - The values are the respective counts of these values in each column

```
1 data = pd.DataFrame({"Qu1": [1, 3, 4, 3, 4],  
2                     "Qu2": [2, 3, 1, 2, 3],  
3                     "Qu3": [1, 5, 2, 4, 4]})  
4  
5 data  
6 Output:  
7      Qu1  Qu2  Qu3  
8  0      1    2    1  
9  1      3    3    5  
10 2     4    1    2  
11 3     3    2    4  
12 4     4    3    4  
13  
14 data["Qu1"].value_counts().sort_index()  
15 Output: Qu1  
16      1    1  
17      3    2  
18      4    2  
19      Name: count, dtype: int64
```

```
1 result = data.apply(pd.value_counts).fillna(0)  
2  
3 result  
4 Output:  
5      Qu1  Qu2  Qu3  
6  1    1.0  1.0  1.0  
7  2    0.0  2.0  1.0  
8  3    2.0  2.0  0.0  
9  4    2.0  0.0  2.0  
10 5    0.0  0.0  1.0
```

Unique Values, Value Counts, and Membership

- There is also a **DataFrame.value_counts** method, but it computes counts considering each row of the DataFrame as a tuple to determine the number of occurrences of each distinct row
 - In this case, the result has an index representing the distinct rows as a hierarchical index

```
1 data = pd.DataFrame({"a": [1, 1, 1, 2, 2], "b": [0, 0, 1, 0, 0]})  
2  
3 data  
4 Output:  
5      a   b  
6    0   0  
7    1   0  
8    2   1  
9    3   0  
10   4   0  
11  
12 data.value_counts()  
13 Output:  
14      a   b  
15    1   0     2  
16    2   0     2  
17    1   1     1  
18 Name: count, dtype: int64
```



The End

THE END!!!

